1    This application is submitted in the name of the following inventor(s):

2

| 3 *Inventor* | *Citizenship* | *Residence City and State* |
|---|---|---|
| 4 Vijayan RAJAN | India | Sunnyvale, California |
| 5 Jeff KIMMEL | United States | San Jose, California |
| 6 Steven R. KLEIMAN | United States | Los Altos, California |
| 7 Joseph CARADONNA | United States | Santa Clara, California |

8

9       The assignee is *Network Appliance, Inc.*, a California corporation having an

10  office at 495 East Java Drive, Sunnyvale, CA 94089.

11

12                    TITLE OF THE INVENTION

13

14  Symmetric Multiprocessor Synchronization Using Migrating Scheduling Domains

15

16                    BACKGROUND OF THE INVENTION

17

18  *1.    Field of the Invention*

19

20       This invention relates to symmetric multiprocessor synchronization and im-

21  plicit synchronization of resources using migrating scheduling domains, as described

22  herein.

1   *2.     Related Art*

2

3          In computer systems having multiple processors with concurrent execution,

4    it is desirable to use as much of the parallelism as possible from the multiple processors.

5    One problem with using the parallelism of multiple processors is that of designing soft-

6    ware to make use of that parallelism. For example, software that was designed for use

7    with a uniprocessor system often does not exploit the parallelism of a multiprocessor

8    system to the fullest extent possible.

9

10         A first known method is to redesign or rewrite software originally designed

11   for use with a uniprocessor system, so as to make use of the advantages of a multiproces-

12   sor system. While this known method does generally achieve the goal of using the ad-

13   vantages of a multiprocessor system, it is subject to several drawbacks. First, it is ex-

14   tremely expensive, in that it uses relatively large amounts of (human design and coding)

15   resources for redesigning or rewriting program code. Second, it is sometimes then neces-

16   sary to maintain two different code bases, one for uniprocessor systems and one for mul-

17   tiprocessor systems, also resulting in additional expense and use of human design and

18   coding resources.

19

20         A second known method is to introduce (into software originally designed

21   for use with a uniprocessor system) those explicit synchronization methods for maintain-

22   ing integrity of resources to be shared among multiple processors. While this known

1    method generally achieves the same goal with relatively less expense and consumption of

2    resources than a complete redesign or rewrite of the software code base, it also suffers

3    from several drawbacks. First, it introduces a relatively large amount of new code subject

4    to possible error in coding. Second, it introduces additional processor and memory usage

5    to implement known explicit synchronization methods (such as locking mechanisms),

6    with resulting slowing of the system using those known explicit synchronization methods.

7    The second drawback is particularly exacerbated for resources that are primarily used by

8    only one software element, but find occasional use by a second software element; the first

9    software element pays the price of known explicit synchronization methods for each use

10   of the resource, even though contention for that resource might be relatively rare. Moreo-

11   ver, these drawbacks for this second known method are also applicable to the first known

12   method, as a new design would likely employ explicit synchronization methods.

13

14   A third known method is to identify (within software originally designed for

15   use with a uniprocessor system) those functional elements that can each independently

16   operate without using known explicit synchronization methods. An example of this third

17   known method is shown in U.S. Patent 5,485,579; in that patent, each separated func-

18   tional element is bound to a single processor in a multiprocessor system, so that the sys-

19   tem can assure that each processor is performing functions that do not require known ex-

20   plicit synchronization methods. While this method generally achieves the goal of using

21   the advantages of a multiprocessor system, it is subject to several drawbacks. First, the

22   mapping between separated functional elements and processors is 1:1, so if the number of

1     separated functional elements differs from the number of processors, the system will ei-

2     ther underutilize at least some of the processors or underperform the functions of at least

3     some of the separated functional elements. Second, there is no provision for load balanc-

4     ing among the multiple processors. Third, there is no useful technique for altering the

5     code base so as to make use of greater parallelism, without resorting to the first known

6     method described above.

7

8     Accordingly, it would be advantageous to provide a technique for schedul-

9     ing a set of tasks in a multiprocessor system that is not subject to drawbacks of the known

10     art. In a preferred embodiment, this is achieved using a method and system for providing

11     parallel execution of those tasks while implicitly synchronizing access to a set of re-

12     sources (such as data structures or hardware devices) used by that system.

13

14     SUMMARY OF THE INVENTION

15

16     The invention provides a method and system for scheduling a set of tasks in

17     an MP (multiprocessor) system, and provides parallel execution of those tasks while im-

18     plicitly synchronizing access to a set of resources (such as data structures or hardware de-

19     vices) used by that system. Tasks in the MP system are each assigned to a scheduling do-

20     main, thus associating those tasks with a set of resources controlled by that domain. A

21     scheduler operating at each processor in the MP system implicitly synchronizes those re-

22     sources controlled by each domain, by scheduling only one task for each domain to exe-

1    cute concurrently in the system. Because each instance of the scheduler selects which task

2    is next run independently of its processor, each domain can migrate from one processor to

3    another; thus, each domain can have a task executing on any processor, so long as no do-

4    main has two tasks executing concurrently in the system. Thus, domains (and their tasks)

5    are not bound to any particular processor. Hence the method and system are symmetric.

6

7         A preferred embodiment uses the implicit synchronization enforced by the

8    scheduler for resources controlled by a single domain, and performs explicit synchroniza-

9    tion only for resources shared by more than one domain. When a resource is needed by a

10    task in a first domain but controlled by a second domain, the task can re-designate itself

11    (and thus switch) from the first to the second domain; this allows execution by other tasks

12    in the first domain, while preserving domain scheduling heuristics. This technique pro-

13    vides for implicit synchronization of resources controlled by the first domain, so that ex-

14    plicit synchronization is not needed.

15

16         A preferred embodiment can designate a set of tasks known to be MP-safe

17    (safe for use in an MP system) to not be assigned to any particular domain, as those tasks

18    can be executed concurrently with all other tasks. MP-safe tasks can include: (a) tasks

19    that do not use any resources controlled by a particular domain, such as tasks that perform

20    only computation and/or keep only their own data structures; (b) tasks that already use

21    explicit synchronization for resources they need; and (c) tasks that are otherwise deter-

22    mined by programmers to not require explicit synchronization.

1    Those of ordinary skill in the art will recognize, after perusal of this appli-

2  cation, the many advantages provided by the invention. These include, but are not limited

3  to, the following:

4

5    • The invention provides for MP-safe operation of applications originally designed

6      for a UP system without having to substantially rewrite those applications.

7

8    • The invention provides for MP-safe operation of applications having distinguish-

9      able functional elements, without having to incur overhead associated with explicit

10     synchronization.

11

12   • The invention provides for incrementally altering program code originally de-

13     signed for a UP (uniprocesor) system to execute in an MP-safe manner on an MP

14     system.

15

16   • The invention provides for increased parallel use of an MP system without having

17     to assign any functional elements to particular processors. This is particularly valu-

18     able where the number of functional elements are greater than the number of par-

19     ticular processors.

20

21   • The invention provides for improved load balancing among the processors in an

22     MP system.

1    The invention has general applicability to applications of any kind execut-

2    ing in an MP system in which the scheduler provides implicit synchronization using do-

3    mains. Although a preferred embodiment is described with regard to a file server, there is

4    no particular limitation of the invention to file servers or similar devices. Techniques used

5    by a preferred embodiment of the invention for implicit synchronization, domain migra-

6    tion, domain switching, and the like can be used in contexts other than the specific appli-

7    cations disclosed herein.

8

9                    BRIEF DESCRIPTION OF THE DRAWINGS

10

11        Figure 1 shows a conceptual diagram of a multiprocessor system including

12    a set of migrating scheduling domains.

13

14        Figure 2 shows a diagram of scheduling domains, the tasks and resources

15    that reside in them, and the resource access methods.

16

17        Figure 3 shows a conceptual diagram of symmetric scheduling using a set

18    of migrating scheduling domains.

19

20        Figure 4 shows a conceptual diagram of symmetric scheduling in which

21    tasks can switch scheduling domains.

22

1        Figure 5 shows a process flow diagram of a symmetric scheduler using a set

2    of migrating scheduling domains.

3

4        Figure 6 shows a conceptual diagram of explicit synchronization and im-

5    plicit synchronization.

6

7              DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

8

9        In the following description, a preferred embodiment of the invention is de-

10    scribed with regard to preferred process steps and data structures. Those skilled in the art

11    would recognize after perusal of this application that embodiments of the invention can

12    be implemented using one or more general purpose processors or special purpose proces-

13    sors or other circuits adapted to particular process steps and data structures described

14    herein, and that implementation of the process steps and data structures described herein

15    would not require undue experimentation or further invention.

16

17    *Related Information*

18

19        Inventions described herein can be used in conjunction with inventions de-

20    scribed in the following application(s):

21

1    • Application Serial No. _____, Express Mail Mailing No.

2    _____, filed the same day, in the name of inventors Christopher

3    PEAK, Sathya BETTADAPURA, and Jeffrey KIMMEL, titled "Automatic Verifi-

4    cation of Scheduling Domain Consistency", attorney docket number 103.1066.01.

5

6    Each of these application(s) is hereby incorporated by reference as if fully

7    set forth herein. They are collectively referred to as the "incorporated disclosures."

8

9    *Lexicography*

10

11    The following terms refer or relate to aspects of the invention as described

12    below. The descriptions of general meanings of these terms are not intended to be limit-

13    ing, only illustrative.

14

15    • **concurrent** — in general, overlapping in time. Concurrent access to a resource in-

16    cludes all forms of access in which the data structure or hardware device might be

17    found in an inconsistent state as a consequence of access by more than one task.

18

19    • **execute** or **run** — in general, when the scheduler grants control of a processor to a

20    task to perform the instructions which makeup its program.

21

- **executable** or **runnable** – when a task is ready to execute, but has not yet been granted control of a processor (due to priority or scheduling domain restrictions).

- **MP-safe** — in general, suitable for operation in a multiprocessor environment without corruption or loss of data, or improper contention for resources. Techniques that serve to prevent improper concurrent access render their subject MP-safe.

- **resource** — in general, any software or hardware object for which concurrent access would need some form of synchronization for safe operation in a multiprocessor system. Examples of resources include, but are not limited to: data structures, disk drives, modems, NIC cards, NVRAM, PCI devices, sound cards, and other peripherals.

- **schedule** — in general, to select a task for execution or running by a processor. There is no particular requirement that scheduling involves either (a) operation for a selected period of time, or (b) operation in the future. A task is "scheduled" when it is made runnable.

- **scheduling domain** — in general, a set of tasks and resources selected by designers or program coders for operation in a multiprocessor system, where it's under-

1     stood that only one task in the set is allowed to run and access the resources at any

2     given time. Multiple tasks from different scheduling domains can run concurrently.

3

4     • **symmetric scheduling**— in general, scheduling tasks on each processor inde-

5        pendent of the identity of the processor. In the invention, the method of scheduling

6        performed on each processor is substantially the same, and no tasks are bound to or

7        special to any particular processor. Thus, a task is just as likely to run on one proc-

8        essor as it is on any other.

9

10     • **synchronize** — in general, to protect a resource from improper concurrent access.

11        Known synchronization techniques include explicit synchronization, in which de-

12        signers or program coders include an explicit synchronization mechanism such as a

13        lock or semaphore. The invention provides implicit synchronization, in which re-

14        sources are synchronized by operation of the scheduling domain restrictions de-

15        scribed herein.

16

17     • **task** — in general, any schedulable entity. A task might include a single process, a

18        single thread, or some other individually schedulable entity.

19

20        As noted above, these descriptions of general meanings of these terms are

21   not intended to be limiting, only illustrative. Other and further applications of the inven-

22   tion, including extensions of these terms and concepts, would be clear to those of ordinary

1     skill in the art after perusing this application. These other and further applications are part

2     of the scope and spirit of the invention, and would be clear to those of ordinary skill in the

3     art, without further invention or undue experimentation.

4

5     *System Elements*

6

7           Figure 1 shows a conceptual diagram of a multiprocessor system including

8     a set of migrating scheduling domains.

9

10           A system 100 includes a plurality of processors 110 and a shared memory

11     120.

12

13           Each processor 110 has access to the shared memory 120 and to both (ex-

14     ecutable) tasks 121 and resources 122 therein. The shared memory includes the tasks 121

15     and the resources 112 to be used by those tasks. Tasks 121 and resources 122 are each as-

16     sociated with a single scheduling domain 123 (with exceptions as described below). Thus,

17     for example, in a preferred embodiment there are three scheduling domains 123 called

18     "Network", "Storage", and "Filesystem", each of which has associated therewith one or

19     more tasks 121 and one or more resources 122. Each processor 110 schedules only those

20     tasks 121 within scheduling domains 123 not already in use by other processors 110, so

21     that each scheduling domain 123 is associated with only one processor 110 at a time. At

1   each moment, each processor 110 executes a task 121 associated with only a single one of

2   the scheduling domains 123.

3

4         As described herein, scheduling domains 123 are not bound or fixed to any

5   particular processor 110 forever, or even for any specific time period. Each processor 110

6   schedules its tasks 121 independently of which processor 110 is performing the schedul-

7   ing (that is, scheduling is symmetric with regard to processors 110), so it is possible for

8   any scheduling domain 123 to be associated to any processor 110 from time to time, sub-

9   ject only to the restriction, as further described below, that each scheduling domain 123

10   can be associated with only one processor 110 at a time. Since tasks 121 within a sched-

11   uling domain 123 can be executed at one time by one processor 110 and at another time

12   by a different processor 110, scheduling domains 123 are said to be able to "migrate"

13   from one processor 110 to another.

14

15         As described herein, resources 122 can include data structures, devices, and

16   other objects for which concurrent access by more than one task 121 would need some

17   form of synchronization for safe operation in a multiprocessor system. Although the de-

18   scription herein primarily refers to resources 122 as data structures, there is no particular

19   limitation of resources 122 to only data structures in a general form of the invention.

20

1    *Scheduling Domain Structure*

2

3    Figure 2 shows a diagram of scheduling domains, the tasks and resources

4    that reside in them, and the resource access methods.

5

6    The system 100 includes a plurality of scheduling domains 123, each in-

7    cluding a set of possibly runnable tasks 121 and a set of resources 122. Resource access is

8    possible via one or more of, or some combination of, the following:

9

10   • A task 121 can directly read data from, or write data to, a resource 122 in the same

11   scheduling domain 123.

12

13   When a task 121 reads data from, or writes data to, a resource 122 in the same

14   scheduling domain 123, that interface between the task 121 and the resource 122 is

15   considered MP-safe (safe for multiprocessor operation), because resources 122 are

16   implicitly synchronized by the system 100, as described below.

17

18   • A task 121 can directly read data from, or write data to, a resource 122 in a differ-

19   ent scheduling domain 123 (or a resource 122 not assigned to any scheduling do-

20   main 123).

21

When a task 121 directly reads data from, or writes data to, a resource 122 in a different scheduling domain 123 (or a resource 122 not assigned to any scheduling domain 123), that interface between the task 121 and the resource 122 is not considered MP-safe, unless the resource 122 is explicitly synchronized using an explicit synchronization technique, as described below.

- A task 121 can access data from another domain by proxy. That is, a message can be sent to a server task, which resides in the target domain. The server task accesses the data directly (on behalf of it's client), and then passes the results back to the client task.

- A task 121 can switch from a first scheduling domain 123 to a second scheduling domain 123.

Figure 4 shows a conceptual diagram of scheduling in which tasks can switch scheduling domains.

Each task 121 includes an associated label indicating to which scheduling domain 123 the task 121 belongs. The task's application code 311 in the application layer 310 performs a "relabel-this-task" system call 313 to scheduler code 322 in the kernel layer 320.

1       Similar to the description with regard to figure 3, the application code 311 makes

2       the system call 312 to transfer control to the scheduler code 322. The scheduler

3       code 322 alters the task's associated label, thus changing the task's scheduling

4       domain 123, and places the transferor task 121 on the runnable queue 550 (in its

5       proper location in response to its relative priority). The relabel code 322 transfers

6       control to the scheduler code 321, which (as described with regard to figure 3) se-

7       lects a next task 121 to execute, performs a context switch into that next task 121,

8       and "returns" to a second set of application code 312 in that next task 121.

9

10      Unlike the operation of scheduler code 321 with regard to figure 3, in this figure

11      the scheduler code 321 "returns" back to the same application code 311 in the

12      "next" task 121. As a consequence, the task 121 continues to execute the same ap-

13      plication code 311, but with a new label indicating that the task 121 belongs to a

14      different scheduling domain 123.

15

16      Since the scheduler refuses to schedule more than one task 121 from the originat-

17      ing scheduling domain 123, if the task 121 attempts to switch domains into a

18      scheduling domain 123 that is already in use (that is, there is already a task exe-

19      cuting that belongs to the target scheduling domain 123), the task 121 will not be

20      selected for execution, and will block until at least that scheduling domain 123 is

21      freed for use. Because the task 121 has been relabeled when the scheduler selects

22      the next task 121, the scheduler is willing to schedule another task 121 from the

"old" scheduling domain 121. For example, if a task 121 in the "Network" scheduling domain 123 switches into the "Storage" scheduling domain 123, the "Network" scheduling domain 123 is freed for use by other tasks 121 in that domain, while the "Storage" scheduling domain 123 is occupied by the switching task 121 and not available for use by other tasks 121.

- A task 121 can release control of its processor allowing another task to run, which is either in the same or a different scheduling domain 123.

When a task 121 releases its processor, the task 121 invokes the scheduler (that is, the part of the operating system that performs task scheduling), which determines which task 121 is next to run on the same processor 110. The scheduler chooses a next task 121 to run independent of which processor 110 it is invoked on, but in response to which scheduling domains are being executed on other processors 110. The scheduler does not select any task 121 that would cause a scheduling domain 123 to be executed on more than one processor 110 at once.

Transfer of control from one task 121 to another can result in a different task from the same domain being run immediately. However, it is possible that the scheduler will select a different, higher priority, task 121 from a different scheduling domain so that the transfer of control will result in that higher priority task 121 being run instead.

1    *Symmetric Scheduling*

2

3    Figure 3 shows a conceptual diagram of symmetric scheduling using a set

4    of migrating scheduling domains.

5

6    Each task 121 is able to invoke a scheduler on the processor 110 on which it

7    is executing. The scheduler is independent of which processor on which it is running (al-

8    though in alternative embodiments, it may be that the scheduler takes into account which

9    processor on which it is running, such as for load-balancing or cache affinity purposes,

10   while still allowing scheduling domains 123 to migrate from one processor 110 to an-

11   other). The figure shows a task 121 invoking the scheduler and causing a context switch

12   to a different task 121 in a different scheduling domain 123.

13

14   The task 121 includes application code 311 in an application layer 310 of a

15   computing system on its processor 110. The application code 311 performs a system call

16   that results in invoking the scheduler code 321 in a kernel layer 320 of the computing

17   system. Application layers, kernel layers, and system calls are known in the art of operat-

18   ing systems.

19

20   The application code 311 makes the system call 312 which transfers control

21   to the scheduler code 321. The scheduler code 321 selects a next task 121 to execute, per-

22   forms a context switch into that next task 121, and "returns" to a second set of application

1    code 312 in that next task 121. Unlike known schedulers, the scheduler code 321 selects

2    only those tasks 121 capable of running without causing a scheduling domain 123 to be

3    running on two different processors 110 concurrently.

4

5    *Method of Scheduling*

6

7         Figure 5 shows a process flow diagram of a symmetric scheduler using a set

8    of migrating scheduling domains.

9

10        A method 500 includes a set of flow points and process steps as described

11   herein.

12

13        At a flow point 510, application code 311 makes the system call 312 to in-

14   voke scheduler code 321, and the scheduler is entered. The scheduler uses a queue 550 of

15   runnable tasks 121, each of which is labeled with a scheduling domain 123 associated

16   with that task 121. The queue 550 includes a head 551 of the queue 550, which identifies

17   a particular task 121.

18

19        At a step 511, the scheduler examines the queue 550 for a next runnable

20   task 121. If there is no such task 121 (that is, the queue 550 is empty or has been com-

21   pletely traced down), the scheduler goes into an idle mode and proceeds with the flow

22   point 510. (Thus, the idle mode can be entered in one of two different ways: first, if there

1    is no next task 121 to run, that is, the runnable queue 550 is empty; second, if there is no

2    task 121 on the runnable queue 550 capable of being run, due to scheduling domain 123

3    restrictions.) If there is such a task 121, the scheduler proceeds with the next step.

4

5    In alternative embodiments, the runnable queue 550 may be separated into a

6    separate runnable queue per scheduling domain 123. This implementation may optimize

7    (speed-up) the scheduler lookup and queue functions.

8

9    At a step 512, the scheduler examines the task 121 at the identified position

10   in the queue 550, and determines which scheduling domain 123 the task 121 is associated

11   with.

12

13   At a step 513, the scheduler determines if that scheduling domain 123 is

14   available for scheduling. If so, the scheduler proceeds with the flow point 514. If not, the

15   scheduler proceeds with the step 511.

16

17   At a step 514, the scheduler prepares to run the selected new task 121. The

18   scheduler performs a context switch into the selected new task 121, and proceeds with the

19   flow point 520.

20

21   At a flow point 520, processor 110 is running the selected task's application

22   code 312.

1 *Synchronization , Explicit or Implicit*

2

3        Figure 6 shows a conceptual diagram of explicit synchronization and im-

4 plicit synchronization.

5

6        With explicit synchronization, a first task 121 and a second task 121 each

7 attempt to access a shared resource 122, such as a data structure. To prevent improper

8 concurrent access to the shared resource 122, each task 121 makes explicit calls 601 to a

9 synchronization mechanism 602. The synchronization mechanism 602 might include a

10 lock, a semaphore, a monitor, or other methods known in the art of operating systems.

11

12        With implicit synchronization, it is assumed by the application that the

13 scheduler will provide the synchronization, by not running multiple tasks in the same do-

14 main concurrently. The first task 121 and the second task 121 each have an associated

15 scheduling domain 123. If the two scheduling domains 123 are different, that indicates a

16 designer's or program coder's decision that the two tasks 121 will not perform improper

17 concurrent access to the shared resource 122 (in alternative embodiments, different

18 scheduling domains 123 may indicate that if there is any improper concurrent access to

19 the shared resource 122, no harm will come to the system 100). If the two scheduling do-

20 mains 123 are the same, that indicates a designer's or program coder's decision that the

21 two tasks 121 might perform improper concurrent access to the shared resource 122, thus

22 that the two tasks 121 are not allowed to execute concurrently.

1        The scheduler prevents concurrent execution of the two tasks 121, and

2   therefore prevents concurrent access to the shared resource 122, as a consequence of the

3   steps 512 and 513 described above. Because the scheduler refuses to schedule two tasks

4   121 for concurrent execution on different processors 110 when those two tasks 121 are

5   associated with the same scheduling domain 123, the two tasks 121 are implicitly syn-

6   chronized with regard to the shared resource 122. The resource 122 is therefore also asso-

7   ciated with the same scheduling domain 123. Lack of improper concurrent access to the

8   resource 122 is therefore an emergent consequence of the scheduler's behavior in refus-

9   ing to concurrently schedule tasks 121 from the same scheduling domain 123.

10

11   *Tasks and Resources Not In Any Domain*

12

13        Tasks 121 or resources 122 can also be declared by the designer or program

14   coder to not be in any scheduling domain 123.

15

16        If a task 121 is not in any domain, the designer or program coder thus indi-

17   cates that the task 121 is MP-safe, that is, that running the task 121 will not result in any

18   improper concurrent access to any resources 122. A task 121 can be declared MP-safe for

19   one or more of the following reasons:

20

21      • The task 121 does not use any resources 122, or all of its resources are internal to

22          the task 121.

1   For example, the task 121 might perform only operations on data structures given

2   to it by a requesting task 121, or the task 121 might use only its own internal data

3   structures (dynamically allocated by the operating system, not static). One example

4   could be when the task 121 performs an operation that does not involve use of re-

5   sources 122, such as a pure calculation.

6

7   • The task 121 already uses explicit synchronization methods to assure it is MP-safe.

8

9   For example, the task 121 might use locks or semaphores to maintain all its data

10  structures or other resources 122 MP-safe. A first example could be when the task

11  121 is part of the kernel and was already designed to be MP-safe. A second exam-

12  ple could be when the task 121 is designed or coded explicitly for safe use in a

13  multiprocessor system.

14

15  • Designers or program coders decide that the task's nature is such that it will not

16  involve improper concurrent access.

17

18  For example, although the task 121 is not formally MP-safe, and does use re-

19  sources 122 that are not formally MP-safe, designers or program coders knowing

20  the nature of the system 100 might have decided that the task 121 will not ever be

21  run, in normal operation, in an unsafe circumstance.

22

1 • Designers or program coders decide that the task's nature is such that, even if it

2 does involve improper concurrent access, that improper concurrent access would

3 be harmless.

4

5 For example, although the task 121 might be able to perform improper concurrent

6 access, designers or program coders knowing the nature of the system 100 might

7 have decided that in that event, the resources 122 that are improperly accessed do

8 not need to be maintained consistent, or that there are adequate recovery proce-

9 dures from that event.

10

11 If a resource 122 is not in any domain, the designer or program coder thus

12 indicates that the resource 122 is MP-safe, that is, that using the resource 122 will not re-

13 sult in any improper concurrent access. A resource 122, and any library code which may

14 maintain the resource, can be declared MP-safe for one or more of the following reasons:

15

16 • The resource 122 already uses explicit synchronization methods to assure it is MP-

17 safe.

18

19 For example, the resource 122 might require locks or semaphores to be accessed.

20

21 • Designers or program coders decide that the resource's nature is such that it will

22 not involve improper concurrent access.

1    For example, although the resource 122 is not formally MP-safe, and might be

2    used concurrently by multiple tasks 121, designers or program coders knowing the

3    nature of the system 100 might have decided that the task 121 will not ever be run,

4    in normal operation, in an unsafe circumstance.

5

6    • Designers or program coders decide that the resource's nature is such that, even if

7    it does involve improper concurrent access, that improper concurrent access would

8    be harmless.

9

10    For example, although the resource 122 might be able to be concurrently accessed,

11    designers or program coders knowing the nature of the system 100 might have de-

12    cided that in that event, the resource 122 does not need to be maintained consis-

13    tent, or that there are adequate recovery procedures from that event.

14

15    *Tasks and Resources In More Than One Domain*

16

17    Tasks 121 or resources 122 can also be declared by the designer or program

18    coder to be in more than one scheduling domain 123.

19

20    In a preferred embodiment, a task 121 can perform a system call to "grab" a

21    second scheduling domain 123 for a period of time. In this case, both the task's first

1 scheduling domain 123 and the task's second scheduling domain 123 are not free for con-

2 current use by other tasks 121.

3

4 In a preferred embodiment, designers or program coders can declare a re-

5 source 122 to be "part of" both a first scheduling domain 123 and a second scheduling

6 domain 123.

7

8 *Preemptive Multitasking*

9

10 A preferred embodiment described herein uses non-preemptive multitask-

11 ing; that is, a task 121 only blocks if it makes the appropriate system call to the scheduler

12 to block itself and allow another task 121 to run.

13

14 In a system 100 using preemptive multitasking, a task 121 can be preempted

15 "against its will," that is, without the task 121 necessarily having a chance to assure that

16 all its resources 122 or other data structures are in order for another task 121 to run. In

17 this case, the task 121 might have left one or more resources 122 in a state that disallows

18 other tasks 121 from the same scheduling domain 123 from accessing those resources

19 122. Accordingly, in alternative embodiments using preemptive multitasking, when a task

20 121 is preempted, it becomes the only next task 121 from its scheduling domain 123 able

21 to next run. Thus, the scheduler will select, from each scheduling domain 123, the pre-

22 empted task 121 over all other tasks 121 in that scheduling domain 123.

1 *Generality of the Invention*

2

3        The invention has general applicability to applications of any kind execut-

4 ing in an MP system in which the scheduler provides implicit synchronization using do-

5 mains. Although a preferred embodiment is described with regard to a file server, there is

6 no particular limitation of the invention to file servers or similar devices. Techniques used

7 by a preferred embodiment of the invention for implicit synchronization, domain migra-

8 tion, domain switching, and the like can be used in contexts other than the specific appli-

9 cations disclosed herein.

10

11        The invention is generally applicable to all applications capable of being

12 run in a multiprocessor system, and to any multiprocessor system in which the scheduler

13 (or equivalent part of an operating system) can be used to enforce implicit synchroniza-

14 tion as described herein.

15

16        Other and further applications of the invention in its most general form

17 would be clear to those skilled in the art after perusal of this application. The invention

18 would be usable for such other and further applications without undue experimentation or

19 further invention.

20

1        Although preferred embodiments are disclosed herein, many variations are

2    possible which remain within the concept, scope and spirit of the invention; these varia-

3    tions would be clear to those skilled in the art after perusal of this application.